

**Preconditioners for ill-conditioned Toeplitz matrices**

Daniel Potts and Gabriele Steidl

230/97

December 1997

# Preconditioners for ill-conditioned Toeplitz matrices

Daniel Potts  
Medizinische Universität zu Lübeck  
Institut für Mathematik  
Wallstr. 40  
D-23560 Lübeck  
potts@informatik.mu-luebeck.de

and

Gabriele Steidl  
Universität Mannheim  
Fakultät für Mathematik und Informatik  
D-68131 Mannheim  
steidl@kiwi.math.uni-mannheim.de

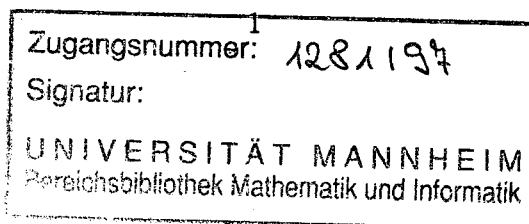
**Abstract.** This paper is concerned with the solution of systems of linear equations  $A_N x = b$ , where  $\{A_N\}_{N \in \mathbb{N}}$  denotes a sequence of positive definite Hermitian ill-conditioned Toeplitz matrices arising from a (real-valued) nonnegative generating function  $f \in C_{2\pi}$  with zeros. We construct positive definite Hermitian preconditioners  $M_N$  such that the eigenvalues of  $M_N^{-1} A_N$  are clustered at 1 and the corresponding PCG-method requires only  $\mathcal{O}(N \log N)$  arithmetical operations.

We sketch how our preconditioning technique can be extended to symmetric Toeplitz systems, doubly symmetric block Toeplitz systems with Toeplitz blocks and non-Hermitian Toeplitz systems.

Numerical tests confirm the theoretical expectations.

1991 *Mathematics Subject Classification.* 65F10, 65F15, 65T10.

*Key words and phrases.* Ill-conditioned Toeplitz matrices, CG-method, clusters of eigenvalues, preconditioners.



# 1 Introduction

Systems of linear equations

$$\mathbf{A}_N \mathbf{x} = \mathbf{b}$$

with positive definite Hermitian Toeplitz matrices  $\mathbf{A}_N$  arise in a variety of applications in mathematics and engineering (see [9] and the references therein). Along with stabilization techniques for direct fast and superfast Toeplitz solvers, preconditioned conjugate gradient methods (PCG-methods) and other iterative methods have attained much attention during the last years. As essential computational effort, the CG-method requires the multiplication of a vector with the matrix  $\mathbf{A}_N$  in each iteration step. For Toeplitz matrices  $\mathbf{A}_N$ , the multiplication with a vector can be computed with  $\mathcal{O}(N \log N)$  arithmetical operations by fast Fourier transforms (FFT). The number of iteration steps of the CG-method depends on the distribution of the eigenvalue of  $\mathbf{A}_N$ . In particular, it holds (see [1], p. 573)

**Theorem 1.1.** Let  $\mathbf{A}_N$  be a positive definite Hermitian  $(N, N)$ -matrix which has  $p$  and  $q$  isolated large and small eigenvalues, respectively:

$$\begin{aligned} 0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_q < a \leq \lambda_{q+1} \leq \dots \leq \lambda_{N-p} \leq b \\ < \lambda_{N-p+1} \leq \lambda_{N-p+2} \leq \dots \leq \lambda_N \quad (0 < a < b < \infty). \end{aligned} \quad (1.1)$$

Let  $[x]$  denote the smallest integer  $\geq x$ . Then the CG-method for the solution of  $\mathbf{A}_N \mathbf{x} = \mathbf{b}$  requires at most

$$n = \left\lceil \left( \ln \frac{2}{\tau} + \sum_{k=1}^q \frac{b}{\lambda_k} \right) / \ln \frac{1 + (\frac{a}{b})^{1/2}}{1 - (\frac{a}{b})^{1/2}} \right\rceil + p + q$$

iteration steps to achieve precision  $\tau$ , i.e.

$$\frac{\|\mathbf{x}_n - \mathbf{x}\|_A}{\|\mathbf{x}_0 - \mathbf{x}\|_A} \leq \tau,$$

where  $\|\mathbf{x}\|_A := \sqrt{\bar{\mathbf{x}}' \mathbf{A}_N \mathbf{x}}$  and where  $\mathbf{x}_n$  denotes the numerical solution after  $n$  iteration steps.

Let  $\{\sigma_k^N\}_{k=1}^N$  be a sequence of real numbers and let  $\gamma_N(\varepsilon)$  denote the number of those among  $\sigma_k^N$  ( $k = 1, \dots, N$ ) which are outside the interval  $(p - \varepsilon, p + \varepsilon)$ . If  $\gamma_N(\varepsilon) < K(\varepsilon)$ , where  $K(\varepsilon)$  is independent of  $N$ , then we say that the values  $\sigma_k^N$  are *clustered at  $p$*  [29]. If the eigenvalues of a sequence of  $(N, N)$ -matrices  $\mathbf{A}_N$  are clustered at 1, then the CG-method converges superlinearly (see [11]).

For a sequence of  $(N, N)$ -Toeplitz matrices  $\mathbf{A}_N = \mathbf{A}_N(f)$  ( $N \in \mathbb{N}$ ) generated by a function  $f \in C_{2\pi}$ , it is well-known that the eigenvalues are distributed as  $f$  [29, 16]. Let

$$f_{\min} := \min\{f(x) : x \in [0, 2\pi)\}, \quad f_{\max} := \max\{f(x) : x \in [0, 2\pi)\}.$$

Then the eigenvalues of  $A_N(f)$  are contained in  $[f_{\min}, f_{\max}]$ . If  $f > 0$ , then by Theorem 1.1 the number of iteration steps of the CG-method is independent of  $N$  and the CG-method requires only  $\mathcal{O}(N \log N)$  arithmetical operations.

The situation changes completely, if we allow  $f \geq 0$  to have zeros. In this case, the CG-method converges very slowly with increasing  $N$ . To accelerate the convergence of the CG-method, several authors proposed preconditioners for Toeplitz systems. Clearly, the multiplication with the preconditioned matrix should also only require  $\mathcal{O}(N \log N)$  arithmetical operations. Therefore, two types of preconditioners were mainly exploited for linear Toeplitz systems, namely so-called "Strang-preconditioners" [14, 27, 12]

$$M_N(S_N f, F_N) := F_N \operatorname{diag} \left( (S_N f) \left( \frac{2\pi j}{N} \right) \right)_{j=0}^{N-1} \bar{F}_N, \quad (1.2)$$

where  $S_N f$  denotes the  $(N-1)$ -th Fourier sum of  $f$  and optimal preconditioners [13]

$$M_N^{\mathcal{O}}(F_N) := F_N \delta(\bar{F}_N A_N F_N) \bar{F}_N, \quad (1.3)$$

where  $\delta(A) := \operatorname{diag}(a_{kk})_{k=0}^{N-1}$  and  $a_{kk}$  are the diagonal entries of  $A$ . Here  $F_N$  denotes the  $N$ -th Fourier matrix

$$F_N := \frac{1}{\sqrt{N}} (e^{-2\pi i j k / N})_{j,k=0}^{N-1}.$$

If  $f > 0$ , then both preconditioners  $M_N$  are positive definite and the eigenvalues of the preconditioned matrices  $M_N^{-1} A_N$  are clustered at 1.

Unfortunately, if  $f$  has zeros, then the eigenvalues of the preconditioned matrices do not fulfil (1.1). Moreover, the Strang-preconditioner may not be positive definite.

Therefore, E. E. Tyrtyshnikov [28] replaced the above preconditioners by improved circulants. Other authors [5, 7, 8, 25] suggested banded Toeplitz matrices or multigrid methods [15] as preconditioners.

In this paper, we propose simple preconditioners which are up to multiplications with unitary diagonal matrices, again circulant matrices. In particular, if  $f(2\pi j/N) > 0$  for all  $j = 0, \dots, N-1$ , then we obtain our preconditioners by replacing  $S_N f$  in (1.2) by  $f$ . In Section 3, we prove that our preconditioners lead to superlinear convergence of the corresponding PCG-method.

Our idea can be extended to (real) symmetric Toeplitz matrices, non-Hermitian Toeplitz matrices and doubly symmetric block Toeplitz matrices with Toeplitz blocks. We sketch the various generalizations in Section 4. Writing this paper, we became aware of the preprint [19] of T. Huckle located at his home page, where the author suggests a trigonometric preconditioner with respect to the discrete sine transform of type I which is similar to our trigonometric preconditioners in Section 4. However, our initial approach in the complex case and our proofs are different from [19].

Numerical tests for Hermitian and symmetric Toeplitz matrices as well as for non-symmetric Toeplitz matrices and doubly symmetric block Toeplitz matrices with Toeplitz blocks in Section 5 demonstrate the quality of our new preconditioners.

## 2 Construction of preconditioners

Let  $C_{2\pi}$  and  $L_{2\pi}^p$  ( $1 \leq p < \infty$ ) denote the Banach spaces of  $2\pi$ -periodic continuous functions and of  $2\pi$ -periodic Lebesgue measurable functions with finite integral  $\int_0^{2\pi} |f(x)|^p dx$ , respectively. By  $\mathbf{o}_N$ , we denote the vector consisting of  $N$  zeros and by  $\mathbf{I}_N$  the  $(N, N)$ -identity matrix.

We are interested in the solution of Hermitian Toeplitz systems

$$\mathbf{A}_N(f) \mathbf{x} = \mathbf{b}, \quad \mathbf{A}_N(f) := (a_{j-k})_{j,k=0}^{N-1}, \quad (2.1)$$

where the sequence  $\{\mathbf{A}_N(f)\}_{N=1}^\infty$  of Toeplitz matrices is generated by a nonnegative function  $f \in C_{2\pi}$ , i.e.

$$a_k = a_k(f) := \frac{1}{2\pi} \int_0^{2\pi} f(x) e^{-ikx} dx.$$

Then it holds for  $\mathbf{u} = (u_j)_{j=0}^{N-1} \in \mathbb{C}^N$  that

$$\begin{aligned} \bar{\mathbf{u}}' \mathbf{A}_N(f) \mathbf{u} &= \sum_{j=0}^{N-1} \sum_{k=0}^{N-1} \bar{u}_j u_k a_{j-k} = \frac{1}{2\pi} \sum_{j=0}^{N-1} \sum_{k=0}^{N-1} \bar{u}_j u_k \int_0^{2\pi} f(x) e^{-i(j-k)x} dx \\ &= \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 f(x) dx \geq 0 \end{aligned} \quad (2.2)$$

such that the Toeplitz matrices  $\mathbf{A}_N(f)$  are positive semidefinite. Moreover, if  $f > 0$  on a set of positive Lebesgue measure, then following lemma states that the matrices  $\mathbf{A}_N(f)$  are positive definite such that (2.1) can be solved by the CG-method.

**Lemma 2.1.** Let  $f \in L_{2\pi}^1$  be a nonnegative function, where the set  $\{x \in [0, 2\pi] : f(x) > 0\}$  has a positive Lebesgue measure. Then the corresponding Toeplitz matrices  $\mathbf{A}_N(f)$  are positive definite.

Lemma 2.1 was proved in [7]. However, the proof is very short such that we include it in this paper.

**Proof.** Let  $N \in \mathbb{N}$  be fixed. By the above considerations, it remains to show that 0 is not eigenvalue of  $\mathbf{A}_N(f)$ . Assume that  $\mathbf{A}_N(f)$  has eigenvalue 0. Then there exists  $\mathbf{u} \in \mathbb{C}^N$   $\mathbf{u} \neq \mathbf{o}_N$  such that

$$\bar{\mathbf{u}}' \mathbf{A}_N(f) \mathbf{u} = \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 f(x) dx = 0.$$

Since the integrand is nonnegative almost everywhere, the integrand must be zero almost everywhere. Consequently,

$$\left| \sum_{k=0}^{N-1} u_k e^{ikx} \right| = 0$$

on the set  $\{x \in [0, 2\pi] : f(x) > 0\}$  of positive Lebesgue measure. But this implies the contradiction  $u = o_N$ .  $\blacksquare$

By Theorem 1.1, the convergence of the CG-method depends on the distribution of the eigenvalues of  $A_N(f)$ . Unfortunately, if the generating  $f \in C_{2\pi}$  has zeros, then the CG-method converges very slowly. To accelerate the convergence of the CG-method we are looking for a suitable preconditioners  $M_N(f)$  of  $A_N(f)$ . Having Theorem 1.1 in mind, we want to find a Hermitian positive definite matrix  $M_N(f)$  such that the number of isolated eigenvalues of  $M_N(f)^{-1}A_N(f)$  is independent of  $N$ . If in addition the multiplication with  $M_N(f)^{-1}$  requires only  $\mathcal{O}(N \log N)$  arithmetical operations, then (2.1) can be solved by the PCG-method with  $\mathcal{O}(N \log N)$  arithmetical operations. For the construction of  $M_N(f)$  we consider (2.2). In the following, we assume that  $f$  has only a finite number of zeros. Then we can choose an equispaced grid

$$x_l := \frac{2\pi l}{N} + w \quad (w \in [0, \frac{2\pi}{N}); l = 0, \dots, N-1)$$

such that

$$f(x_l) > 0 \quad (l = 0, \dots, N-1). \quad (2.3)$$

Approximating the integral on the right-hand side of (2.2) by the trapezoidal rule with respect to the above grid, we obtain

$$\begin{aligned} \bar{u}' A_N(f) u &= \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 f(x) dx \\ &\approx \frac{1}{N} \sum_{l=0}^{N-1} \left| \sum_{k=0}^{N-1} u_k e^{ikx_l} \right|^2 f(x_l) \\ &= \sum_{l=0}^{N-1} f(x_l) \frac{1}{\sqrt{N}} \left( \sum_{j=0}^{N-1} \bar{u}_j e^{-2\pi i l j / N} e^{-i j w} \right) \frac{1}{\sqrt{N}} \left( \sum_{k=0}^{N-1} u_k e^{2\pi i k l / N} e^{i k w} \right) \quad (2.4) \\ &= (\mathbf{F}_N \mathbf{W}_N \bar{\mathbf{u}})' \mathbf{D}_N \bar{\mathbf{F}}_N \bar{\mathbf{W}}_N \mathbf{u} \\ &= \bar{\mathbf{u}}' \mathbf{M}_N(f) \mathbf{u} \end{aligned}$$

with the diagonal matrices

$$\mathbf{W}_N := \text{diag} (e^{-ikw})_{k=0}^{N-1}, \quad \mathbf{D}_N := \text{diag} (f(x_l))_{l=0}^{N-1}$$

and with

$$\mathbf{M}_N(f) = \mathbf{M}_N(f, \mathbf{F}_N) := \mathbf{W}_N \mathbf{F}_N \mathbf{D}_N \bar{\mathbf{F}}_N \bar{\mathbf{W}}_N. \quad (2.5)$$

By (2.3), the matrix  $\mathbf{M}_N(f)$  is Hermitian and positive definite. Setting  $\mathbf{v} := \mathbf{M}_N(f)^{1/2} \mathbf{u}$ , we get

$$\bar{\mathbf{v}}' \mathbf{M}_N(f)^{-1/2} A_N(f) \mathbf{M}_N(f)^{-1/2} \mathbf{v} \approx \bar{\mathbf{v}}' \mathbf{v}$$

such that by properties of the Rayleigh quotient,  $\mathbf{M}_N(f)$  seems to be a good preconditioner of  $\mathbf{A}_N(f)$ . Indeed, using FFT, the multiplication with

$$\mathbf{M}_N(f)^{-1} = \mathbf{W}_N \mathbf{F}_N \mathbf{D}_N^{-1} \bar{\mathbf{F}}_N \bar{\mathbf{W}}_N$$

takes only  $\mathcal{O}(N \log N)$  arithmetical operations. In the next section, we prove that the eigenvalues of  $\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f)$  are clustered at 1.

We mention that our preconditioner  $\mathbf{M}_N(f)$  is closely related to the Strang-preconditioner  $\mathbf{M}_N(\mathcal{S}_N f) = \mathbf{M}_N(\mathcal{S}_N f, \mathbf{F}_N)$  in (1.2). By orthogonality of the functions  $e^{ijx}$  ( $j \in \mathbb{Z}$ ) in  $L^2_{2\pi}$ , it is easy to check that (2.2) can be replaced by

$$\bar{\mathbf{u}}' \mathbf{A}_N(f) \mathbf{u} = \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 (\mathcal{S}_N f)(x) dx$$

with

$$(\mathcal{S}_N f)(x) := \sum_{j=-(N-1)}^{N-1} a_j e^{ijx}.$$

Now the above quadrature formula (2.4) with  $w = 0$  and with  $\mathcal{S}_N f$  instead of  $f$  leads to the Strang-preconditioner. Clearly, if  $f$  is a trigonometric polynomial of degree  $< N$  and if  $f(2\pi l/N) > 0$  ( $l = 0, \dots, N-1$ ), then  $\mathbf{M}_N(\mathcal{S}_N f) = \mathbf{M}_N(f)$ .

However, for arbitrary nonnegative functions  $f \in C_{2\pi}$ , the matrix  $\mathbf{M}_N(\mathcal{S}_N f)$  may be not positive definite. This is one reason for the introduction of  $\mathbf{M}_N(f)$ .

### 3 Clustering of the eigenvalues of $\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f)$

We rewrite (2.4) as

$$\begin{aligned} \bar{\mathbf{u}}' \mathbf{A}_N(f) \mathbf{u} &= \sum_{j=0}^{N-1} \sum_{k=0}^{N-1} \bar{u}_j u_k a_{j-k} \\ &\approx \sum_{j=0}^{N-1} \sum_{k=0}^{N-1} \bar{u}_j u_k \tilde{a}_{j-k} = \bar{\mathbf{u}}' \mathbf{M}_N(f) \mathbf{u} \end{aligned} \quad (3.1)$$

with

$$\tilde{a}_k = \tilde{a}_k(f) := \frac{1}{N} \sum_{l=0}^{N-1} f(x_l) e^{-2\pi i l k / N} e^{-ikw}$$

and ask for the approximation error. Assume that  $f \in C_{2\pi}$  is a function of bounded

variation. Replacing  $f(x_l)$  by the Fourier series of  $f$  at  $x_l$ , we obtain

$$\begin{aligned}
\tilde{a}_k &= \frac{1}{N} \sum_{l=0}^{N-1} \sum_{j \in \mathbb{Z}} a_j e^{ijx_l} e^{-2\pi i l k / N} e^{-ikw} \\
&= \sum_{j=0}^{N-1} a_j e^{-i w k} e^{i w j} \left( \frac{1}{N} \sum_{l=0}^{N-1} e^{-2\pi i l k / N} e^{2\pi i l j / N} \right) \\
&\quad + \sum_{j=0}^{N-1} \sum_{r \in \mathbb{Z} \setminus \{0\}} a_{j+rN} e^{-i w k} e^{i w (j+rN)} \left( \frac{1}{N} \sum_{l=0}^{N-1} e^{-2\pi i l k / N} e^{2\pi i l j / N} \right) \\
&= a_k + \sum_{r \in \mathbb{Z} \setminus \{0\}} a_{k+rN} e^{i w r N}.
\end{aligned} \tag{3.2}$$

This is the well-known *aliasing effect*. Set

$$b_k = b_k(f) := \sum_{r \in \mathbb{Z} \setminus \{0\}} a_{k+rN}(f) e^{i w r N}. \tag{3.3}$$

Then it follows by (3.2) that

$$\mathbf{A}_N(f) = \mathbf{M}_N(f) + \mathbf{B}_N(f), \quad \mathbf{B}_N(f) := -(b_{j-k})_{j,k=0}^{N-1}. \tag{3.4}$$

Thus

$$\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f) = \mathbf{I}_N + \mathbf{M}_N(f)^{-1} \mathbf{B}_N(f). \tag{3.5}$$

Note that

$$b_k(\mathcal{S}_N f) = \begin{cases} a_{k-N}(f) & k = 1, \dots, N-1, \\ a_{k+N}(f) & k = -1, \dots, 1-N, \\ 0 & k = 0, \end{cases}$$

which describes the approximation error in case of the Strang-preconditioner.

**Lemma 3.1.** Let  $p_s$  be a nonnegative trigonometric polynomial of degree  $\leq s$ , where  $2s \leq N$ . Then at most  $2s$  eigenvalues of  $\mathbf{M}_N(p_s)^{-1} \mathbf{A}_N(p_s)$  differ from 1.

**Proof:** By (3.3), it follows that  $b_k = 0$  for  $|k| \leq N-1-s$ . Consequently,  $\mathbf{B}_N(f)$  has rank  $2s$ . Now the assertion follows by (3.5).  $\blacksquare$

For the proof of our main theorem we need the following

**Lemma 3.2.** Let  $g \in C_{2\pi}$  be nonnegative functions, where the set  $\{x \in [0, 2\pi] : f(x) > 0\}$  has a positive Lebesgue measure. Furthermore let  $h \in C_{2\pi}$  be a positive function with  $h_{\min} > 0$  and let  $f := gh$ . Then, for any  $N \in \mathbb{N}$ , the eigenvalues of  $\mathbf{A}_N(g)^{-1} \mathbf{A}_N(f)$  lie in the interval  $[h_{\min}, h_{\max}]$ .



Lemma 3.2 was proved for example in [5]. We want to give a different simple proof based on the theorem of mean.

**Proof:** Applying the theorem of mean in

$$\bar{u}' A_N(f) u = \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 f(x) dx,$$

we obtain that

$$\bar{u}' A_N(f) u = h_* \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{k=0}^{N-1} u_k e^{ikx} \right|^2 g(x) dx$$

with  $h_* \in [h_{\min}, h_{\max}]$ . This can be rewritten as

$$\bar{u}' A_N(f) u = h_* \bar{u}' A_N(g) u.$$

By Lemma 2.1, the matrix  $A_N(g)$  is positive definite such that for  $u \neq o_N$

$$h_* = \frac{\bar{u}' A_N(f) u}{\bar{u}' A_N(g) u}.$$

By properties of the Rayleigh quotient this yields the assertion. ■

In the following, we restrict our attention to nonnegative functions  $f \in C_{2\pi}$  having a zero of *even order*  $2s$  ( $s \in \mathbb{N}$ ) in  $x = 0$ . The clustering of the eigenvalues of  $M_N(f)^{-1} A_N(f)$  for arbitrary functions

$$f(x) = (x - x_1)^{2s_1} \dots (x - x_m)^{2s_m} \tilde{f}(x), \quad (\tilde{f} > 0)$$

follows in a similar way.

With  $f \in C_{2\pi}$  or better with the order  $2s$  of the zero  $x = 0$  of  $f$ , we associate the nonnegative trigonometric polynomial

$$p_s(x) := (2 - 2 \cos x)^s = (2 - e^{ix} - e^{-ix})^s \quad (3.6)$$

of degree  $s$  which has a zero of the same order  $2s$  in  $x = 0$ .

Now we can prove our main result.

**Theorem 3.3.** Let  $f \in C_{2\pi}$  be a nonnegative function with a zero of order  $2s$  ( $s \in \mathbb{N}$ ) in  $x = 0$ . Let  $A_N(f)$  denote the corresponding Toeplitz matrices with preconditioners  $M_N(f)$  defined by (2.5). Then the matrices  $M_N(f)^{-1} A_N(f)$  have the following properties:

- i) The eigenvalues of  $M_N(f)^{-1} A_N(f)$  are clustered at 1.
- ii) Let  $p_s$  denote the associated trigonometric polynomial (3.6) of  $f$  and let  $h := f/p_s$ . Then, for  $N \geq 2s$ , at most  $2s$  eigenvalues of  $M_N(f)^{-1} A_N(f)$  are not contained in the

interval  $[\frac{h_{\min}}{h_{\max}}, \frac{h_{\max}}{h_{\min}}]$ .

**Proof:** In this proof, we denote by  $R_N(m)$  arbitrary  $(N, N)$ -matrix of rank  $m$ .

1. To prove ii), we use the decomposition [4]

$$\frac{\bar{u}' A_N(f) u}{\bar{u}' M_N(f) u} = \frac{\bar{u}' A_N(f) u}{\bar{u}' A_N(p_s) u} \frac{\bar{u}' A_N(p_s) u}{\bar{u}' M_N(p_s) u} \frac{\bar{u}' M_N(p_s) u}{\bar{u}' M_N(f) u} \quad (u \neq o_N). \quad (3.7)$$

By (3.4) and Lemma 3.1, it holds that

$$A_N(p_s) = M_N(p_s) + R_N(2s).$$

The Hermitian matrix  $R_N(2s)$  can be written as

$$R_N(2s) = \bar{Q}' D_+ Q + \bar{Q}' D_- Q = R_N(s_1) + R_N(2s - s_1) \quad (3.8)$$

where  $Q$  denotes a unitary matrix and where  $D_+$ ,  $D_-$  are diagonal matrices containing the  $s_1$  positive eigenvalues and the  $2s - s_1$  negative eigenvalues of  $R_N(2s)$  as diagonal entries, respectively. Setting

$$a(u) := \frac{\bar{u}' A_N(f) u}{\bar{u}' A_N(p_s) u}, \quad m(u) := \frac{\bar{u}' M_N(p_s) u}{\bar{u}' M_N(f) u},$$

we obtain by (3.7) and (3.8) that

$$\frac{\bar{u}' A_N(f) u}{\bar{u}' M_N(f) u} = a(u) m(u) + a(u) \frac{\bar{u}' R_N(s_1) u}{\bar{u}' M_N(f) u} + a(u) \frac{\bar{u}' R_N(2s - s_1) u}{\bar{u}' M_N(f) u}.$$

By Lemma 3.2 and by construction of  $M_N$ , it follows for all  $u \in \mathbb{C}^N$  ( $u \neq o_N$ ) that

$$a(u) \in [h_{\min}, h_{\max}], \quad m(u) \in [\frac{1}{h_{\max}}, \frac{1}{h_{\min}}].$$

Since for all  $u \in \mathbb{C}^N$  ( $u \neq o_N$ )

$$\frac{\bar{u}' R_N(s_1) u}{\bar{u}' M_N(f) u} \geq 0, \quad \frac{\bar{u}' R_N(2s - s_1) u}{\bar{u}' M_N(f) u} \leq 0,$$

we get

$$\frac{\bar{u}' A_N(f) u}{\bar{u}' M_N(f) u} \leq \frac{h_{\max}}{h_{\min}} + \frac{\bar{u}' [h_{\max} R_N(s_1) + h_{\min} R_N(2s - s_1)] u}{\bar{u}' M_N(f) u}$$

and further

$$\frac{\bar{u}' [A_N(f) - (h_{\max} R_N(s_1) + h_{\min} R_N(2s - s_1))] u}{\bar{u}' M_N(f) u} \leq \frac{h_{\max}}{h_{\min}}.$$

The matrix  $h_{\max} R_N(s_1) + h_{\min} R_N(2s - s_1)$  has again  $s_1$  positive and  $2s - s_1$  negative eigenvalues. Thus, by properties of the Rayleigh quotient and by Weyl's theorem [17], p. 184, at most  $s_1$  eigenvalues of  $M_N(f)^{-1} A_N(f)$  are larger than  $\frac{h_{\max}}{h_{\min}}$ .

Similarly, we obtain that

$$\frac{\bar{u}' [A_N(f) - (h_{\min} R_N(s_1) + h_{\max} R_N(2s - s_1))] u}{\bar{u}' M_N(f) u} \geq \frac{h_{\min}}{h_{\max}}.$$

Thus, at most  $2s - s_1$  eigenvalues of  $M_N(f)^{-1} A_N(f)$  are smaller than  $\frac{h_{\min}}{h_{\max}}$ . Consequently, at most  $2s$  eigenvalues of  $M_N(f)^{-1} A_N(f)$  are not contained in  $[\frac{h_{\min}}{h_{\max}}, \frac{h_{\max}}{h_{\min}}]$ .

2. By definition,  $h = f/p_s$  is a continuous positive function. Since the trigonometric polynomials are dense in  $C_{2\pi}$ , for all  $\varepsilon > 0$ , there exist a positive trigonometric polynomial  $q$  of degree  $n = n(\varepsilon)$  such that

$$q(x) - \frac{1}{2} \varepsilon h_{\min} \leq h(x) \leq q(x) + \frac{1}{2} \varepsilon h_{\min} \quad (3.9)$$

for all  $x \in [0, 2\pi)$ . Thus, since  $p_s \geq 0$ ,

$$q p_s - \frac{1}{2} \varepsilon h_{\min} p_s \leq f \leq q p_s + \frac{1}{2} \varepsilon h_{\min} p_s. \quad (3.10)$$

Regarding (2.2), we obtain by the inequality of the right-hand side

$$\bar{u}' A_N(f) u \leq \bar{u}' A_N(q p_s) u + \frac{1}{2} \varepsilon h_{\min} \bar{u}' A_N(p_s) u,$$

and further, since  $M_N(f)$  is positive definite, for all  $u \in \mathbb{C}^N$  ( $u \neq o_N$ )

$$\frac{\bar{u}' A_N(f) u}{\bar{u}' M_N(f) u} \leq \frac{\bar{u}' A_N(q p_s) u}{\bar{u}' M_N(f) u} + \frac{1}{2} \varepsilon h_{\min} \frac{\bar{u}' A_N(p_s) u}{\bar{u}' M_N(f) u}. \quad (3.11)$$

Now it holds by (3.4) and Lemma 3.1 that

$$A_N(p_s) = M_N(p_s) + R_N(2s). \quad (3.12)$$

Moreover, we have by [3] that

$$A_N(q p_s) = A_N(q) A_N(p_s) + R_N(2n + 2s).$$

By (3.12), this can be written as

$$\begin{aligned} A_N(q p_s) &= (M_N(q) + R_N(2n)) (M_N(p_s) + R_N(2s)) + R_N(2n + 2s) \\ &= M_N(q) M_N(p_s) + R_N(m) \end{aligned} \quad (3.13)$$

with a Hermitian matrix  $R_N(m)$  of rank  $m \leq 4n + 4s + \min\{2n, 2s\}$ . Substituting (3.12) and (3.13) in (3.11), we obtain

$$\begin{aligned} \frac{\bar{u}' A_N(f) u}{\bar{u}' M_N(f) u} &\leq \frac{\bar{u}' M_N(q) M_N(p_s) u}{\bar{u}' M_N(f) u} + \frac{\bar{u}' R_N(m) u}{\bar{u}' M_N(f) u} \\ &\quad + \frac{1}{2} \varepsilon h_{\min} \frac{\bar{u}' M_N(p_s) u}{\bar{u}' M_N(f) u} + \frac{1}{2} \varepsilon h_{\min} \frac{\bar{u}' R_N(2s) u}{\bar{u}' M_N(f) u} \end{aligned}$$

and since

$$\frac{\bar{u}' M_N(p_s) u}{\bar{u}' M_N(f) u} \leq \frac{1}{h_{\min}}$$

further

$$\frac{\bar{u}' [A_N(f) - R_N(\tilde{m})] u}{\bar{u}' M_N(f) u} \leq \frac{\bar{u}' M_N(q) M_N(p_s) u}{\bar{u}' M_N(f) u} + \frac{1}{2} \varepsilon$$

with  $\tilde{m} \leq m+2s$ . Setting  $v := M_N(p_s)^{1/2} u$  and using that  $M_N(f) = M_N(h) M_N(p_s)$ , we get

$$\frac{\bar{u}' [A_N(f) - R_N(\tilde{m})] u}{\bar{u}' M_N(f) u} \leq \frac{\bar{v}' M_N(q) v}{\bar{v}' M_N(h) v} + \frac{1}{2} \varepsilon. \quad (3.14)$$

Finally, we have by (3.9) and by definition of  $M_N$ , for all  $v \in \mathbb{C}^N$  ( $v \neq o_N$ ) that

$$\bar{v}' M_N(q) v \leq \bar{v}' M_N(h) v + \frac{1}{2} \varepsilon h_{\min} \bar{v}' v$$

and further since  $0 < \frac{\bar{v}' v}{\bar{v}' M_N(h) v} \leq \frac{1}{h_{\min}}$  that

$$\frac{\bar{v}' M_N(q) v}{\bar{v}' M_N(h) v} \leq 1 + \frac{1}{2} \varepsilon.$$

Using the above inequality in (3.14), we obtain

$$\frac{\bar{u}' [A_N(f) - R_N(\tilde{m})] u}{\bar{u}' M_N(f) u} \leq 1 + \varepsilon.$$

Similarly, we conclude from the left-hand inequality of (3.10) that

$$\frac{\bar{u}' [A_N(f) - R_N(\tilde{m})] u}{\bar{u}' M_N(f) u} \geq 1 - \varepsilon.$$

Consequently, at most  $\tilde{m}$  eigenvalues of  $M_N(f)^{-1} A_N(f)$  are not contained in  $[1 - \varepsilon, 1 + \varepsilon]$ . This completes the proof. ■

By Theorem 3.3, Theorem 1.1 and construction of  $M_N(f)$  in (2.5), our PCG-method converges superlinearly and requires only  $\mathcal{O}(N \log N)$  arithmetical operations.

**Remark:** Unfortunately, we cannot find a similar proof for nonnegative functions  $f \in C_{2\pi}$  having not only zeros of even order. The reason therefore is that there does not exist a nonnegative trigonometric polynomial which has a zero of odd order in  $x = 0$ . Consequently, we cannot produce an equivalent of (3.6). Our numerical tests show that our preconditioners work well also in the odd case. However, for the matrices  $A_N(f)$  generated by the function

$$f(x) = \sqrt{2 - 2 \cos x} = |2 \sin \frac{x}{2}|,$$

the number  $n$  of eigenvalues of  $M_N^{-1}(f) A_N(f)$  which are not contained in the interval  $(1 - \varepsilon, 1 + \varepsilon)$  grows as follows:

$N$	32	64	128	256	512
$\varepsilon = 10^{-3}$	7	8	9	10	11
$\varepsilon = 10^{-5}$	10	12	13	15	17

At first glance it seems that the eigenvalues of  $M_N(f)^{-1}A_N(f)$  are not clustered at 1.  $\square$

## 4 Generalizations of the preconditioning technique

In this section, we sketch how our preconditioners can be generalized to the following settings:

- $A_N(f)$  are (real) symmetric Toeplitz matrices,
- $A_N(f)$  are non-Hermitian Toeplitz matrices,
- $A_{M,N}(f)$  are doubly symmetric block Toeplitz matrices with Toeplitz blocks.

First, we suppose in addition to Section 2 that the generating function  $f \in C_{2\pi}$  of the matrices  $A_N(f)$  is even. Then

$$a_k = a_k(f) = \frac{2}{\pi} \int_0^\pi f(x) \cos kx \, dx$$

and the Toeplitz matrices  $A_N(f) \in \mathbb{R}^{N,N}$  are symmetric. In this case, the multiplication of a vector with  $A_N(f)$  can be realized using *fast trigonometric transforms* instead of fast Fourier transforms. See the remark after Lemma 4.1. In this way, in the iterative solution of (2.1), complex arithmetic can be completely avoided. This is one of the reasons to look for preconditioners of type (2.5), where the Fourier matrix  $F_N$  is replaced by trigonometric matrices corresponding to fast trigonometric transforms.

In practice, four discrete sine transforms (DST) and four discrete cosine transforms (DCT) were applied (see [30]). Any of these eight trigonometric transforms can be realized with  $\mathcal{O}(N \log N)$  arithmetical operations (see for example [2], [26]). Likewise, we can define preconditioners with respect to any of these transforms. Here we refer to the extensive examinations in [22]. In this paper, we restrict our attention to the so-called DST-II and DCT-II, which are determined by the following transform matrices:

$$\begin{aligned} \text{DCT-II} \quad : \quad C_N^{II} &:= \left(\frac{2}{N}\right)^{1/2} \left( \varepsilon_j^N \cos \frac{j(2k+1)\pi}{2N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N}, \\ \text{DST-II} \quad : \quad S_N^{II} &:= \left(\frac{2}{N}\right)^{1/2} \left( \varepsilon_{j+1}^N \sin \frac{(j+1)(2k+1)\pi}{2N} \right)_{j,k=0}^{N-1} \in \mathbb{R}^{N,N}, \end{aligned}$$

where  $\varepsilon_k^N := 2^{-1/2}$  ( $k = 0, N$ ) and  $\varepsilon_k^N := 1$  ( $k = 1, \dots, N-1$ ). Moreover, we use the DCT-I with transform matrix

$$\tilde{C}_{N+1}^I := \left( (\varepsilon_k^N)^2 \cos \frac{jk\pi}{N} \right)_{j,k=0}^N.$$

The matrices  $\mathbf{C}_N^{II}$  and  $\mathbf{S}_N^{II}$  are orthogonal and

$$\tilde{\mathbf{C}}_{N+1}^I \tilde{\mathbf{C}}_{N+1}^I = \frac{N}{2} \mathbf{I}_{N+1}. \quad (4.1)$$

The eight trigonometric transforms are closely related to Toeplitz matrices [23]. In particular, it holds for the DCT-II and the DST-II:

**Lemma 4.1.** Let  $\text{stoep } \mathbf{a}'$  and  $\text{shank } \mathbf{a}'$  denote a symmetric Toeplitz matrix and a per-symmetric Hankel matrix with first row  $\mathbf{a}'$ , respectively. Then there exist the following relations between trigonometric transforms and symmetric Toeplitz matrices:

$$\begin{aligned} (\mathbf{C}_N^{II})' \mathbf{D} \mathbf{C}_N^{II} &= \frac{1}{2} \text{stoep}(a_0, \dots, a_{N-1}) + \frac{1}{2} \text{shank}(a_1, \dots, a_{N-1}, 0), \\ (\mathbf{S}_N^{II})' \tilde{\mathbf{D}} \mathbf{S}_N^{II} &= \frac{1}{2} \text{stoep}(a_0, \dots, a_{N-1}) - \frac{1}{2} \text{shank}(a_1, \dots, a_{N-1}, 0) \end{aligned}$$

with

$$\begin{aligned} \mathbf{D} &:= \text{diag}(d_0, \dots, d_{N-1})', \quad \tilde{\mathbf{D}} := \text{diag}(d_1, \dots, d_N)', \\ \mathbf{d} &= (d_0, \dots, d_N)' := \tilde{\mathbf{C}}_N^I (a_0, \dots, a_{N-1}, 0)'. \end{aligned}$$

**Proof:** We restrict the proof to the DCT-II. To simplify the notation, we drop the index  $N$  and set  $\mathbf{C} := \mathbf{Z}_N \mathbf{C}_N^{II}$ , and  $\mathbf{D} := \text{diag } \mathbf{d}$  with

$$\mathbf{Z}'_N := (\mathbf{I}_N, \mathbf{o}_N) \in \mathbb{R}^{N, N+1}.$$

Then, by

$$\cos \alpha \cos \beta = \frac{1}{2} \cos(\alpha - \beta) + \frac{1}{2} \cos(\alpha + \beta),$$

the  $(u, v)$ -entry of the matrix  $\mathbf{C}' \mathbf{D} \mathbf{C}$  is

$$(\mathbf{C}' \mathbf{D} \mathbf{C})_{u,v} = \frac{1}{2} \frac{2}{N} \sum_{k=0}^{N-1} (\varepsilon_k^N)^2 d_k \cos \frac{(u-v)k\pi}{N} + \frac{1}{2} \frac{2}{N} \sum_{k=0}^{N-1} (\varepsilon_k^N)^2 d_k \cos \frac{(u+v+1)k\pi}{N},$$

or equivalently, since  $-(-1)^{u-v} d_N = (-1)^{u+v+1} d_N$  for arbitrary  $d_N \in \mathbb{R}$ ,

$$(\mathbf{C}' \mathbf{D} \mathbf{C})_{u,v} = \frac{1}{2} \frac{2}{N} \sum_{k=0}^N (\varepsilon_k^N)^2 d_k \cos \frac{(u-v)k\pi}{N} + \frac{1}{2} \frac{2}{N} \sum_{k=0}^N (\varepsilon_k^N)^2 d_k \cos \frac{(u+v+1)k\pi}{N}.$$

Choosing  $d_N \in \mathbb{R}$  such that

$$\sum_{k=0}^N (\varepsilon_k^N)^2 d_k (-1)^k = 0,$$

we get by symmetry properties of cosine function that

$$\mathbf{C}' \mathbf{D} \mathbf{C} = \frac{1}{2} \text{stoep}(a_0, \dots, a_{N-1}) + \frac{1}{2} \text{shank}(a_1, \dots, a_{N-1}, 0),$$

where

$$(a_0, \dots, a_{N-1}, 0)' = \frac{2}{N} \tilde{C}_{N+1}^I \mathbf{d},$$

i.e. by (4.1),

$$\mathbf{d} = \tilde{C}_{N+1}^I (a_0, \dots, a_{N-1}, 0)'.$$

**Remark:** By Lemma 4.1, it follows that

$$\text{stoep}(a_0, \dots, a_{N-1}) = (C_N^{II})' D C_N^{II} + (S_N^{II})' \tilde{D} S_N^{II}.$$

Thus, if the vector  $\mathbf{d}$  is precomputed by the DCT-I, then the multiplication of a vector with a symmetric Toeplitz matrix of size  $(N, N)$  requires two DCT-II, two DST-II and  $2N$  real multiplications and can therefore be realized in  $\mathcal{O}(N \log N)$  arithmetical operations (see also [23]).  $\square$

Since for even  $f \in C_{2\pi}$ , the  $(N-1)$ -th Fourier sum can be written as

$$(S_N f)(x) = 2 \sum_{k=0}^{N-1} (\varepsilon_k^N)^2 a_k \cos(kx),$$

we obtain by Lemma 4.1 that

$$\begin{aligned} A_N(f) &= (C_N^{II})' (2D) C_N^{II} - \text{shank}(a_1, \dots, a_{N-1}, 0) \\ &= (C_N^{II})' \text{diag} \left( (S_N f)\left(\frac{j\pi}{N}\right) \right)_{j=0}^{N-1} C_N^{II} - \text{shank}(a_1, \dots, a_{N-1}, 0), \quad (4.2) \\ A_N(f) &= (S_N^{II})' (2\tilde{D}) S_N^{II} + \text{shank}(a_1, \dots, a_{N-1}, 0) \\ &= (S_N^{II})' \text{diag} \left( (S_N f)\left(\frac{j\pi}{N}\right) \right)_{j=1}^N S_N^{II} + \text{shank}(a_1, \dots, a_{N-1}, 0). \quad (4.3) \end{aligned}$$

Consequently, we introduce the Strang-type-preconditioners by [24]:

$$\begin{aligned} \text{DCT - II: } M_N(S_N f, C_N^{II}) &:= (C_N^{II})' \text{diag} \left( (S_N f)\left(\frac{j\pi}{N}\right) \right)_{j=0}^{N-1} C_N^{II}, \\ \text{DST - II: } M_N(S_N f, S_N^{II}) &:= (S_N^{II})' \text{diag} \left( (S_N f)\left(\frac{j\pi}{N}\right) \right)_{j=1}^N S_N^{II}. \end{aligned} \quad (4.4)$$

See also [20]. Again, if  $f$  has zeros, then it can not be assured that the Strang-type-preconditioners are positive definite. Therefore, we define similar to (2.5) the preconditioners

$$\begin{aligned} \text{DCT - II: } M_N(f, C_N^{II}) &:= (C_N^{II})' \text{diag} \left( f\left(\frac{j\pi}{N}\right) \right)_{j=0}^{N-1} C_N^{II}, \\ \text{DST - II: } M_N(f, S_N^{II}) &:= (S_N^{II})' \text{diag} \left( f\left(\frac{j\pi}{N}\right) \right)_{j=1}^N S_N^{II}. \end{aligned} \quad (4.5)$$

If  $f(j\pi/N) > 0$  for all  $j = 0, \dots, N-1$ , then  $\mathbf{M}_N(f, \mathbf{C}_N^{II})$  is positive definite. If  $f(j\pi/N) > 0$  for all  $j = 1, \dots, N$ , then  $\mathbf{M}_N(f, \mathbf{S}_N^{II})$  is positive definite.

Note that independent of our results, T. Huckle [19] suggested a preconditioner of type (4.5) with respect to the DST-I.

Clearly, if  $f$  is a trigonometric polynomial of degree  $< N$ , then the Strang-type-preconditioners (4.4) coincide with our preconditioners (4.5). Moreover, we have by (4.2) and (4.3) for trigonometric polynomials  $f = p$  of degree  $\leq s$  ( $2s \leq N$ ) that

$$\mathbf{A}_N(p) = \mathbf{M}_N(p, \mathbf{C}_N^{II}) - \mathbf{R}_N(2s) = \mathbf{M}_N(p, \mathbf{S}_N^{II}) + \mathbf{R}_N(2s).$$

Thus, we can prove in a completely similar way as in Section 3 the following

**Theorem 4.2.** Let  $f \in C_{2\pi}$  be an even nonnegative function with a zero of order  $2s$  ( $s \in \mathbb{N}$ ) in  $x = 0$ . Let  $\mathbf{A}_N(f)$  denote the corresponding Toeplitz matrices with preconditioners  $\mathbf{M}_N(f) = \mathbf{M}_N(f, \mathbf{S}_N^{II})$  defined by (4.4). Then the matrices  $\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f)$  have the following properties:

- i) The eigenvalues of  $\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f)$  are clustered at 1.
- ii) Let  $p_s$  denote the associated trigonometric polynomial (3.6) of  $f$  and let  $h := f/p_s$ . Then, for  $N \geq 2s$ , at most  $2s$  eigenvalues of  $\mathbf{M}_N(f)^{-1} \mathbf{A}_N(f)$  are not contained in the interval  $[\frac{h_{\min}}{h_{\max}}, \frac{h_{\max}}{h_{\min}}]$ .

The PCG-method with our preconditioners can be realized in a more efficient way than the PCG-method with banded Toeplitz matrices as preconditioners:

**Remark:** Our PCG-method requires only two DCT-II, two DST-II and  $\mathcal{O}(N)$  real multiplications in each iteration step. This can be seen for the preconditioner  $\mathbf{M}_N(f, \mathbf{C}_N^{II})$  as follows: Instead of

$$(\mathbf{C}_N^{II})' \mathbf{E}^{-1} \mathbf{C}_N^{II} \left( (\mathbf{C}_N^{II})' \mathbf{D} \mathbf{C}_N^{II} + (\mathbf{S}_N^{II})' \tilde{\mathbf{D}} \mathbf{S}_N^{II} \right) \mathbf{x} = (\mathbf{C}_N^{II})' \mathbf{E}^{-1} \mathbf{C}_N^{II} \mathbf{b}$$

with  $\mathbf{E} := \text{diag} \left( f\left(\frac{j\pi}{N}\right) \right)_{j=0}^{N-1}$ , we solve

$$\mathbf{E}^{-1} \left( \mathbf{D} + \mathbf{C}_N^{II} (\mathbf{S}_N^{II})' \tilde{\mathbf{D}} \mathbf{S}_N^{II} (\mathbf{C}_N^{II})' \right) \tilde{\mathbf{x}} = \tilde{\mathbf{b}}$$

with  $\tilde{\mathbf{x}} := \mathbf{C}_N^{II} \mathbf{x}$  and  $\tilde{\mathbf{b}} := \mathbf{E}^{-1} \mathbf{C}_N^{II} \mathbf{b}$ . The vectors  $\mathbf{d}$ ,  $\tilde{\mathbf{b}}$  and  $\mathbf{x}$  can be precomputed and postcomputed, respectively. See also [18, 19].  $\square$

Next, we are interested in the solution of systems of linear equations  $\mathbf{A}_N(f) \mathbf{x} = \mathbf{b}$  with regular, but non-Hermitian Toeplitz matrices  $\mathbf{A}_N(f)$ . We intend to solve the normal equation

$$\bar{\mathbf{A}}'_N(f) \mathbf{A}_N(f) \mathbf{x} = \bar{\mathbf{A}}'_N(f) \mathbf{b} \quad (4.6)$$

using the PCG-method. By [3], it holds that

$$\bar{\mathbf{A}}'_N(f) \mathbf{A}_N(f) = \mathbf{A}_N(|f|^2) + \mathbf{R}_N + \mathbf{U}_N,$$



with a low rank matrix  $R_N$  and a matrix  $U_N$  of small spectral norm. If  $f = p$  is a trigonometric polynomial of degree  $\leq s$  ( $2s \leq N$ ), then

$$\bar{A}'_N(f)A_N(f) = A_N(|f|^2) + R_N(2s).$$

Assume that  $|f| \in C_{2\pi}$  has only a finite number of zeros. Then we define our preconditioners by

$$M_N(|f|^2, F_N) := W_N F_N \text{diag} \left( |f(\frac{2\pi j}{N} + w)|^2 \right)_{j=0}^{N-1} \bar{F}_N \bar{W}_N \quad (w \in [0, 2\pi/N))$$

if  $A_N(|f|^2)$  is Hermitian and  $|f(\frac{2\pi j}{N} + w)| > 0$  for all  $j = 0, \dots, N-1$  and by

$$\begin{aligned} M_N(|f|^2, C_N^{II}) &:= (C_N^{II})' \text{diag} \left( |f(\frac{\pi j}{N})|^2 \right)_{j=0}^{N-1} C_N^{II}, \\ M_N(|f|^2, S_N^{II}) &:= (S_N^{II})' \text{diag} \left( |f(\frac{\pi j}{N})|^2 \right)_{j=1}^N S_N^{II}, \end{aligned} \quad (4.7)$$

if  $A_N(|f|^2)$  is symmetric and  $|f(\frac{2\pi j}{N})| > 0$  for all  $j = 0, \dots, N-1$  and for all  $j = 1, \dots, N$ , respectively.

Finally, the generalization of our results to doubly symmetric block-Toeplitz systems with Toeplitz blocks is straightforward. We consider systems of linear equations

$$A_{M,N}x = b,$$

where  $A_{M,N}$  denotes a positive definite doubly symmetric block-Toeplitz matrix with Toeplitz blocks (BTTB matrix), i.e.

$$A_{M,N} := (A_{r-s})_{r,s=0}^{M-1} \quad \text{with} \quad A_r := (a_{r,j-k})_{j,k=0}^{N-1}$$

and  $a_{r,j} = a_{|r|,|j|}$ . We assume that the matrices  $A_{M,N}$  are generated by a real-valued  $2\pi$ -periodic continuous even function in two variables, i.e.

$$a_{j,k} := \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^{2\pi} \varphi(s,t) e^{-i(sj+tk)} ds dt.$$

Lemma 4.1 can be extended to BTTB matrices as follows:

$$\begin{aligned} A_{M,N} &= (C_M^{II} \otimes C_N^{II})' D_1 (C_M^{II} \otimes C_N^{II}) + (S_M^{II} \otimes C_N^{II})' D_2 (S_M^{II} \otimes C_N^{II}) \\ &+ (C_M^{II} \otimes S_N^{II})' D_3 (C_M^{II} \otimes S_N^{II}) + (S_M^{II} \otimes S_N^{II})' D_4 (S_M^{II} \otimes S_N^{II}) \end{aligned}$$

with

$$\begin{aligned} D_1 &:= \text{diag} \left( \text{col}(\tilde{a}_{r,j})_{j,r=0}^{N-1, M-1} \right), & D_2 &:= \text{diag} \left( \text{col}(\tilde{a}_{r,j})_{j=0, r=1}^{N-1, M} \right), \\ D_3 &:= \text{diag} \left( \text{col}(\tilde{a}_{r,j})_{j=1, r=0}^{N, M-1} \right), & D_4 &:= \text{diag} \left( \text{col}(\tilde{a}_{r,j})_{j=1, r=1}^{N, M} \right), \end{aligned}$$

$$(\tilde{a}_{r,j})_{j,r=0}^{N,M} := \tilde{C}_{M+1}^I \left( (a_{r,j})_{j,r=0}^{N,M} \right) (\tilde{C}_{N+1}^I)' ,$$

$a_{r,N} := 0$  ( $r = 0, \dots, M$ ) and  $a_{M,j} := 0$  ( $j = 0, \dots, N$ ). Here  $\text{col}: \mathbb{R}^{N,M} \rightarrow \mathbb{R}^{MN}$  is defined by

$$\text{col} (x_{j,k})_{j=0,k=0}^{N-1,M-1} := (x_r)_{r=0}^{MN-1} \quad \text{with} \quad x_{kN+j} := x_{j,k} .$$

Consequently, the multiplication of a vector with a BTTB matrix requires only  $\mathcal{O}(MN \log(MN))$  arithmetical operations. For details see [24]. We define our so-called "level-2" preconditioners by

$$\begin{aligned} \mathbf{M}_N(\varphi, \mathbf{C}_M^{II} \otimes \mathbf{C}_N^{II}) &:= (\mathbf{C}_M^{II} \otimes \mathbf{C}_N^{II})' \text{diag}(\text{col} \left( \varphi\left(\frac{r\pi}{M}, \frac{j\pi}{N}\right) \right)_{j,k=0}^{N-1,M-1}) (\mathbf{C}_M^{II} \otimes \mathbf{C}_N^{II}) , \\ \mathbf{M}_N(\varphi, \mathbf{S}_M^{II} \otimes \mathbf{S}_N^{II}) &:= (\mathbf{S}_M^{II} \otimes \mathbf{S}_N^{II})' \text{diag}(\text{col} \left( \varphi\left(\frac{r\pi}{M}, \frac{j\pi}{N}\right) \right)_{j,k=1}^{N,M}) (\mathbf{S}_M^{II} \otimes \mathbf{S}_N^{II}) . \end{aligned}$$

Using the same arguments as in the remark after Theorem 4.2, we see that our PCG-method requires per iteration step only  $MN$  multiplications more than the conventional CG-method.

## 5 Numerical Examples

In this section, we show the efficiency of our new preconditioning technique by various numerical examples. The fast computation of the preconditioners and the PCG-method were implemented in MATLAB, where the C-programs for the fast Fourier transform and the fast trigonometric transforms were included by cmex. The algorithms were tested on a Sun SPARCstation 20.

As transform length we choose  $N = 2^n$  and as right-hand side  $\mathbf{b}$  of (2.1) the vector consisting of  $N$  entries "1". The PCG-method started with the zero vector and stopped if  $\|\mathbf{r}^{(j)}\|_2 / \|\mathbf{r}^{(0)}\|_2 < 10^{-7}$ , where  $\mathbf{r}^{(j)}$  denotes the residual vector after  $j$  iterations.

We begin with Hermitian ill-conditioned Toeplitz matrices  $\mathbf{A}_N(f)$  arising from the generating function

$$\text{i) } f(x) = (x/2 - \pi/4)^4 \quad (x \in [0, 2\pi)) .$$

The second column of Table 1 shows the number of iterations of the CG-method without preconditioning. The columns 3 and 4 contain the numbers of iterations of the PCG-method with the optimal preconditioner  $\mathbf{M}_N^\mathcal{O}(f, \mathbf{F}_N)$  given by (1.3) and with our preconditioner  $\mathbf{M}_N(f, \mathbf{F}_N)$  defined by (2.5) with  $w := \pi/N$ , respectively.

Next, we consider symmetric Toeplitz matrices  $\mathbf{A}_N$ . We compare the Strang-type-preconditioners (4.4), our preconditioners (2.5) and (4.5) and the optimal trigonometric preconditioners defined by

$$\begin{aligned} \text{DCT-II:} \quad \mathbf{M}_N^\mathcal{O}(\mathbf{C}_N^{II}) &:= (\mathbf{C}_N^{II})' \delta(\mathbf{C}_N^{II} \mathbf{A}_N (\mathbf{C}_N^{II})') \mathbf{C}_N^{II} , \\ \text{DST-II:} \quad \mathbf{M}_N^\mathcal{O}(\mathbf{S}_N^{II}) &:= (\mathbf{S}_N^{II})' \delta(\mathbf{S}_N^{II} \mathbf{A}_N (\mathbf{S}_N^{II})') \mathbf{S}_N^{II} . \end{aligned}$$

$n$	$I_N$	$M_N^{\mathcal{O}}(F_N)$	$M_N(f, F_N)$
4	26	17	11
5	85	36	13
6	349	67	17
7	1570	154	22
8	> 3000	377	26
9	> 3000	995	35
10	> 3000	2220	46

Table 1:  $f(x) = (x/2 - \pi/4)^4 \quad (x \in [0, 2\pi))$

See for example [6, 10, 23]. Our test matrices correspond to the following generating functions:

ii) (see [25]):  $f(x) := (x^2 - 1)^2 \quad (x \in [-\pi, \pi))$ .

In (2.5), we set  $w := \pi/N$ .

iii) (see [25, 7, 8]):  $f(x) := x^4 \quad (x \in [-\pi, \pi))$ .

In (2.5), we set  $w := \pi/N$ .

The Tables 2 and 3 present the number of iteration steps for different preconditioners. The asterix emphasizes that the corresponding preconditioners are not positive definite. Our new preconditioners lead to the best results. Compare also with [25, 7, 8]. Note that by the remark after Theorem 4.2, our PCG-method requires per iteration step only few arithmetical operations more than the conventional CG-method.

Our next test is related to non-symmetric Toeplitz systems. As generating function of  $A_N(f)$  we choose

iv)  $f(x) = x^2 e^{ix} \quad (x \in [-\pi, \pi))$ .

Then, the matrices  $A_N(f)$  have real entries such that we restrict our attention to trigonometric preconditioners. Table 5 compares the PCG-method applied to the normal equation (4.6) with

- the optimal preconditioner of  $A'_N(f)A_N(f)$

$$M_N^{\mathcal{O}_1} := O'_N \delta(O_N A'_N(f) A_N(f) O'_N) O_N,$$

- the optimal preconditioner of  $A_N(|f|^2)$

$$M_N^{\mathcal{O}_2} := O'_N \delta(O_N A_N(|f|^2) O'_N) O_N,$$

- the Strang-type-preconditioner  $M_N(\mathcal{S}_N(|f|^2), O_N)$   
and our preconditioner  $M_N(|f|^2, S_N^{II})$  defined by (4.7).

		$M_N(S_N f, O_N)$		$M_N^O(O_N)$		$M_N(f, O_N)$	
$n$	$I_N$	$C_N^{II}$	$S_N^{II}$	$C_N^{II}$	$S_N^{II}$	$S_N^{II}$	$F_N$
5	25	9*	8*	17	10	5	5
6	69	9*	8*	21	11	5	6
7	190	10*	10*	26	14	7	7
8	457	10*	10*	33	16	8	8
9	> 1000	11	9	43	19	9	9
10	> 1000	10*	10*	59	24	7	7

Table 2:  $f(x) = (x^2 - 1)^2$  ( $x \in [-\pi, \pi)$ )

		$M_N(S_N f, O_N)$		$M_N^O(O_N)$		$M_N(f, O_N)$	
$n$	$I_N$	$C_N^{II}$	$S_N^{II}$	$C_N^{II}$	$S_N^{II}$	$S_N^{II}$	$F_N$
5	33	12*	10*	18	10	6	6
6	116	18*	15*	30	13	7	6
7	487	27*	21*	54	16	8	8
8	>1000	40*	33*	155	19	9	11
9	>1000	115*	63*	376	25	9	13
10	>1000	218*	165*	>1000	32	10	15

Table 3:  $f(x) = x^4$  ( $x \in [-\pi, \pi)$ )

Finally, let us turn to BTTB matrices  $A_{N,N}$ . In our two examples, the matrices  $A_{N,N}$  are generated by the functions

v) (see [21])  $\varphi(s, t) = s^2 t^4$  and  $\psi(s, t) = (s^2 + t^2)^2$  ( $s, t \in [-\pi, \pi)$ ).

Both matrices are ill-conditioned and the CG-method without preconditioning, with Strang-type-preconditioning or with optimal trigonometric preconditioning converges very slowly (see [21, 24]). Our preconditioning determined by (4.7) leads to the number of iterations in Table 5. Again, our PCG-method requires per iteration step only few arithmetical operations more than the conventional CG-method.

$n$	$I_N$	$M_N^{\mathcal{O}_1}(C_N^{II})$	$M_N^{\mathcal{O}_1}(S_N^{II})$	$M_N^{\mathcal{O}_2}(C_N^{II})$	$M_N^{\mathcal{O}_2}(S_N^{II})$	$M_N(\mathcal{S}_N( f ^2), C_N^{II})$	$M_N(\mathcal{S}_N( f ^2), S_N^{II})$	$M_N( f ^2, S_N^{II})$
5	84	29	21	34	18	22*	19*	11
6	311	52	26	64	22	32*	26*	11
7	1226	116	33	139	27	56*	44*	14
8	5220	256	40	324	39	96*	76*	16
9	>10000	664	74	865	55	200*	157*	19
10	>10000	1758	101	2546	78	466*	357*	21

Table 4:  $f(x) := x^2 e^{ix}$  ( $x \in [-\pi, \pi)$ )

$N$	$M_N(\varphi, S_N^{II} \otimes S_N^{II})$	$M_N(\psi, S_N^{II} \otimes S_N^{II})$
8	13	9
16	16	12
32	22	14
64	29	19
128	36	25
256	43	35
512	52	49

Table 5:  $\varphi(s, t) = s^2 t^4$  and  $\psi(s, t) = (s^2 + t^2)^2$  ( $s, t \in [-\pi, \pi]$ )

## References

- [1] O. Axelsson. *Iterative Solution Methods*. Cambridge University Press, Cambridge, 1996.
- [2] G. Baszenski and M. Tasche. Fast polynomial multiplication and convolution related to the discrete cosine transform. *Linear Algebra Appl.*, 252:1 – 25, 1997.
- [3] F. D. Benedetto. Iterative solution of Toeplitz systems by preconditioning with the discrete sine transform. In SPIE 2563, San Diego, 1995.
- [4] F. D. Benedetto and S. S. Capizzano. A unifying approach to abstract matrix algebra preconditioning. *Preprint*, 1997.
- [5] F. D. Benedetto, G. Fiorentino, and S. Serra. C.G. preconditioning for Toeplitz matrices. *Comp. Math. Appl.*, 25:35 – 45, 1993.
- [6] E. Boman and I. Koltracht. Fast transform based preconditioners for Toeplitz equations. *SIAM J. Matrix Anal. Appl.*, 16:628 – 645, 1995.
- [7] R. H. Chan. Toeplitz preconditioners for Toeplitz systems with nonnegative generating functions. *IMA J. Numer. Anal.*, 11:333 – 345, 1991.
- [8] R. H. Chan and K.-P. Ng. Toeplitz preconditioners for hermitian Toeplitz systems. *Linear Algebra Appl.*, 190:181 – 208, 1993.
- [9] R. H. Chan and M. K. Ng. Conjugate gradient methods of Toeplitz systems. *SIAM Review*, 38:427 – 482, 1996.
- [10] R. H. Chan, M. K. Ng, and C. K. Wong. Sine transform based preconditioners for symmetric Toeplitz systems. *Linear Algebra Appl.*, 232:237 – 259, 1996.

- [11] R. H. Chan and G. Strang. Toeplitz systems by conjugate gradients with circulant preconditioner. *SIAM J. Sci. Statist. Comput.*, 10:104 – 119, 1989.
- [12] R. H. Chan and M.-C. Yeung. Circulant preconditioners constructed from kernels. *SIAM J. Numer. Anal.*, 29:1093 – 1103, 1992.
- [13] T. F. Chan. An optimal circulant preconditioner for Toeplitz systems. *SIAM J. Sci. Statist. Comput.*, 9:766 – 771, 1988.
- [14] R.H. Chang and G. Strang. Toeplitz systems by conjugate gradients with circulant preconditioner. *SIAM J. Sci. Statist. Comput.*, 10: 104 – 119, 1989.
- [15] G. Fiorentino and S. Serra. Multigrid methods for Toeplitz matrices. *Calcolo*, 28: 283 – 305, 1992.
- [16] U. Grenander and G. Szegö. *Toeplitz Forms and Their Applications*. University of California Press, Los Angeles, 1958.
- [17] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985.
- [18] T. Huckle. Iterative methods for Toeplitz-like matrices. Report SCCM-94-05, Stanford University, 1994.
- [19] T. Huckle. Iterative methods for ill-conditioned Toeplitz matrices. *Preprint*, 1997.
- [20] T. Kailath and V. Olshevsky. Displacement structure approach to discrete-trigonometric-transform based preconditioners of G. Strang type and of T. Chan type. *Preprint*, 1996.
- [21] M. K. Ng. Band preconditioners for block-Toeplitz-Toeplitz-block-systems. *Linear Algebra Appl.*, 259:307 – 327, 1997.
- [22] D. Potts. *Schnelle Polynomtransformation und Vorkonditionierer für Toeplitz-Matrizen*. PhD thesis, Univ. Rostock, 1998.
- [23] D. Potts and G. Steidl. Optimal trigonometric preconditioners for nonsymmetric Toeplitz systems. *Preprint*, 1996.
- [24] D. Potts, G. Steidl, and M. Tasche. Trigonometric preconditioners for block Toeplitz systems. In *Multivariate Approximation and Splines*, G. Nürnberger, J. W. Schmidt, and G. Walz, (eds), Birkhäuser, Basel, 1997, 219 – 234.
- [25] S. Serra. Optimal, quasi-optimal and superlinear band-Toeplitz preconditioners for asymptotically ill-conditioned positive definite Toeplitz systems. *Math. Comp.*, 66:651 – 665, 1997.
- [26] G. Steidl and M. Tasche. A polynomial approach to fast algorithms for discrete Fourier-cosine and Fourier-sine transforms. *Math. Comp.*, 56:281 – 296, 1991.

- [27] G. Strang. A proposal for Toeplitz matrix calculations. *Studies in Appl. Math.*, 74:171 – 176, 1986.
- [28] E. E. Tyrtyshnikov. Circulant preconditioners with unbounded inverses. *Linear Algebra Appl*, 216:1 – 23, 1995.
- [29] E. E. Tyrtyshnikov. A unifying approach to some old and new theorems on distribution and clustering. *Linear Algebra Appl*, 232:1 – 43, 1996.
- [30] Z. Wang. Fast algorithms for the discrete W transform and for the discrete Fourier transform. *IEEE Trans. Acoust. Speech Signal Process*, 32:803 – 816, 1984.